

A Case for Malleable Thread-Level Linear Algebra Libraries: The LU Factorization with Partial Pivoting

Sandra Catalán^a, José R. Herrero^b, Enrique S. Quintana-Ortí^a, Rafael
Rodríguez-Sánchez^a, Robert van de Geijn^c

^a*Depto. Ingeniería y Ciencia de Computadores, Universidad Jaume I, Castellón, Spain.*

^b*Dept. d'Arquitectura de Computadors, Universitat Politècnica de Catalunya, Spain.*

^c*Department of Computer Science, Institute for Computational Engineering and
Sciences, The University of Texas at Austin, USA.*

Abstract

We propose two novel techniques for overcoming load-imbalance encountered when implementing so-called look-ahead mechanisms in relevant dense matrix factorizations for the solution of linear systems. Both techniques target the scenario where two thread teams are created/activated during the factorization, with each team in charge of performing an independent task/branch of execution. The first technique promotes *worker sharing* (WS) between the two tasks, allowing the threads of the task that completes first to be reallocated for use by the costlier task. The second technique allows a fast task to alert the slower task of completion, enforcing the *early termination* (ET) of the second task, and a smooth transition of the factorization procedure into the next iteration.

The two mechanisms are instantiated via a new *malleable* thread-level implementation of the *Basic Linear Algebra Subprograms* (BLAS), and their benefits are illustrated via an implementation of the LU factorization with partial pivoting enhanced with look-ahead. Concretely, our experimental results on a six core Intel-Xeon processor show the benefits of combining WS+ET, reporting competitive performance in comparison with a task-parallel runtime-based solution.

Email addresses: catalans@uji.es (Sandra Catalán), josepr@ac.upc.edu (José R. Herrero), quintana@uji.es (Enrique S. Quintana-Ortí), rarodrig@uji.es (Rafael Rodríguez-Sánchez), rvdg@cs.utexas.edu (Robert van de Geijn)

Keywords:

1. Introduction

In the 1970s and 80s, the scientific community recognized the value of defining standard interfaces for dense linear algebra (DLA) operations with the introduction of the *Basic Linear Algebra Subprograms* (BLAS) [1, 2, 3]. Ever since, the key to performance portability in this domain has been the development of highly-optimized, architecture-specific implementations of the BLAS, either by hardware vendors (e.g., Intel MKL [4], AMD ACML [5], IBM ESSL [6], and NVIDIA CUBLAS [7]) or independent developers (e.g., GotoBLAS [8, 9], OpenBLAS [10], ATLAS [11], and BLIS [12]).

Multi-threaded instances of the BLAS for current multi-core processor architectures take advantage of the simple data dependencies featured by these operations to exploit loop/data-parallelism at the block level (hereafter referred to as *block-data parallelism*, or BDP). For more complex DLA operations, like those supported by LAPACK [13] and `libflame` [14], exploiting *task-parallelism* with dependencies¹ (TP) is especially efficient when performed by a runtime that semi-automatically decomposes the computation into tasks and orchestrates their dependency-aware scheduling [15, 16, 17, 18]. For the BLAS kernels though, exploiting BDP is still the preferred choice, because it allows tighter control on the data movements across the memory hierarchy and avoids the overhead of a runtime that is unnecessary due to the (mostly) nonexistent data dependencies in the BLAS kernels.

Exploiting both BDP and TP, in a sort of nested parallelism can yield more efficient solutions as the number of cores in processor architectures continues to grow. For example, consider an application composed of two independent tasks, T_A and T_B , both of which are inherently parallel and preceded/followed by synchronization points. In this scenario, exploiting TP only is inefficient, because it can keep at most 2 threads busy. To address this, we could take advantage of the BDP inside T_A and T_B via TP but, given their inherent parallelism, this is likely to incur certain overhead compared

¹We view TP as the type of concurrency present in an operation that can be decomposed into a collection of suboperations (tasks) connected by a rich set of data dependencies. Compared with this, BDP is present when the operation basically consists of a number of independent (inherently parallel) suboperations, each acting on a disjoint block of data.

with a direct exploitation of BDP. Finally, extracting BDP only is surely possible, but it may be less scalable than a nested TP+BDP solution that splits the threads into two teams, and puts them to work on T_A and T_B concurrently.

Let us consider now that the previous scenario occurs during the execution of a complex DLA operation where both tasks, T_A and T_B , can be computed via simple calls to a multi-threaded instance of BLAS. Although the solution seems obvious, exploiting nested TP+BDP parallelism here can still be suboptimal. In particular, *all multi-threaded instances of BLAS offer a rigid interface to control the threading execution of a routine, which only allows one to set the number of threads that will participate before the routine is invoked. Importantly, this number cannot be changed during its execution.* Thus, in case T_A is completed earlier, the team of threads in charge of its execution will remain idle waiting for the completion of T_B , producing a suboptimal execution from a performance perspective.

The scenario that we have described is far from being rare in DLA. To illustrate this, we will employ the LU factorization with partial pivoting for the solution of linear systems [19]. High-performance algorithms for this decomposition consist of a loop-body that processes the matrix from its top-left corner to the bottom-right one, at each iteration computing a panel factorization and updating a trailing submatrix via calls to BLAS. We will review this factorization as a prototypical example to make the following contributions in our paper:

- *Malleable DLA libraries*: We introduce a malleable thread-level implementation of BLIS [12] that allows the number of threads that participate in the execution of a BLAS kernel to dynamically change at execution time.
- *Worker Sharing (WS)*: In case the panel factorization is less expensive than the update of the trailing submatrix, we leverage the malleable instance of the BLAS to improve workload balancing and performance, by allowing the thread team in charge of the panel factorization to be reallocated to the execution of the trailing update.
- *Early Termination (ET)*: To tackle the opposite case, where panel factorization is more expensive than the update of the trailing submatrix, we design an ET mechanism that allows the thread team in charge of the trailing update to communicate the alternative team of this event.

This alert forces an ET of the panel factorization, and the advance of the factorization into the next iteration.

- We perform a comprehensive experimental evaluation on a 6-core Intel Xeon E5-2603 v3 processor, using execution traces to illustrate actual benefits of our approach, and comparing its performance to those obtained with a runtime-based solution using OmpSs [15].

The key to our approach is that we depart from conventional instances of BLAS to instead view the cores/threads as a pool of computational resources that, upon completing the execution of a BLAS/LAPACK routine, can be tapped to participate in the execution of another BLAS/LAPACK routine that is already in progress. This WS supports a dynamic choice of the algorithmic block size as the operation progresses. Furthermore, the same idea carries over to all other major matrix decompositions for the solution of linear systems, such as the QR, Cholesky and LDL^T factorizations [19].

2. The BLIS Implementation of Basic Linear Algebra Kernels

BLIS is a framework that allows developers to rapidly deploy new high-performance implementations of BLAS and BLAS-like operations on current and future architectures [12]. A key property of the BLIS open source effort is that it exposes the internal implementation of the BLAS kernels at a finer-grain level than OpenBLAS or commercial libraries while offering performance that is competitive with GotoBLAS, OpenBLAS, Intel MKL, and ATLAS [20, 21]. We start by reviewing the design principles that underlie BLIS, using the implementation of GEMM as a particular case study.

Consider the matrices $A \in \mathbb{R}^{m \times k}$, $B \in \mathbb{R}^{k \times n}$, and $C \in \mathbb{R}^{m \times n}$. BLIS mimics GotoBLAS to implement the GEMM kernel²

$$C += A \cdot B \tag{1}$$

as three nested loops around a macro-kernel plus two packing routines (see Loops 1–3 in Figure 1). The macro-kernel is then implemented in terms of

²Actually, the kernel in the BLAS interface/BLIS implementation for GEMM computes $C = \alpha C + \beta op(A) \cdot op(B)$, where α, β are scalars, $op(\cdot)$ performs an optional transposition/Hermitian-conjugation, and $op(A)$ is $m \times k$, $op(B)$ is $k \times n$, C is $m \times n$. For simplicity, in the description we address the case where $\alpha = \beta = 1$ and the operator $op(\cdot)$ does not perform any transformation on the input matrix.

```

Loop 1  for  $j_c = 0, \dots, n-1$  in steps of  $n_c$ 
Loop 2    for  $p_c = 0, \dots, k-1$  in steps of  $k_c$ 
           $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$  // Pack into  $B_c$ 
Loop 3    for  $i_c = 0, \dots, m-1$  in steps of  $m_c$ 
           $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$  // Pack into  $A_c$ 
Loop 4    for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$  // Macro-kernel
Loop 5    for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$  // Micro-kernel
           $+= A_c(i_r : i_r + m_r - 1, 0 : k_c - 1)$ 
           $\cdot B_c(0 : k_c - 1, j_r : j_r + n_r - 1)$ 
          

---


        endfor
      endfor
    endfor
  endfor
endfor

```

Figure 1: High performance implementation of GEMM in BLIS. In the code, $C_c \equiv C(i_c : i_c + m_c - 1, j_c : j_c + n_c - 1)$ is just a notation artifact, introduced to ease the presentation of the algorithm. In contrast, A_c, B_c correspond to actual buffers that are involved in data copies.

two additional loops around a *micro-kernel* (Loops 4 and 5 in that figure). The loop ordering embedded in BLIS, together with the packing routines and an appropriate choice of the BLIS cache configuration parameters (n_c , k_c , m_c , n_r and m_r), orchestrate a regular pattern of data transfers across the levels of the memory hierarchy, and amortize the cost of these transfers with enough computation from within the micro-kernel [12] to attain near-peak performance. In most architectures, m_r, n_r are in the range 4–16; m_c, k_c are in the order of a few hundreds; and n_c can be up to a few thousands [12, 20].

The parallelization of BLIS’s GEMM for multi-threaded architectures has been analyzed for conventional symmetric multicore processors [20], modern many-threaded architectures [21], and asymmetric multicore processors [22]. In all these cases, the parallel implementation exploits the BDP exposed by the nested five-loop organization of GEMM, at one or more levels (i.e., loops), using OpenMP or POSIX threads.

A convenient option in most single-socket systems is to parallelize either Loop 3 (indexed by i_c), Loop 4 (indexed by j_r), or a combination of both [20, 21, 22]. For example, when Loop 3 is parallelized, each thread packs a different macro-panel A_c into the L2 cache and executes a different instance of the macro-kernel. In contrast, when Loop 4 is parallelized, different threads will operate on independent instances of the micro-kernel, but access the same macro-panel A_c in the L2 cache.

Consider, for example, a version of BLIS GEMM that extracts BDP from

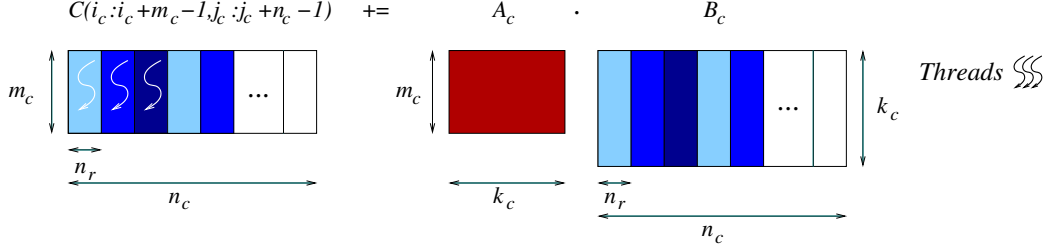


Figure 2: Distribution of the workload among $t = 3$ threads when Loop 4 of BLIS GEMM is parallelized. Different colors in the output C distinguish the panels of this matrix that are computed by each thread as the product of A_c and corresponding panels of the input B_c .

Loop 4 only, to be executed on a multicore architecture with t (physical) cores and one thread mapped per core. The iteration space of Loop 4 is then statically distributed among the t threads in a round-robin fashion, equivalent to the effect attained by adding an OpenMP directive `#pragma omp parallel for`, with a `static` schedule, around that loop; see Figure 2. To improve performance, the packing is also performed in parallel so that, at each iteration of Loop 3, all t threads collaborate to copy and re-organize the entries of $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1)$ into the buffer A_c .

3. Algorithms for the LU Factorization on Multi-threaded Architectures

Given a matrix $A \in \mathbb{R}^{m \times n}$, its LU factorization produces lower and upper triangular factors, $L \in \mathbb{R}^{m \times n}$ and $U \in \mathbb{R}^{n \times n}$ respectively, such that $PA = LU$, where $P \in \mathbb{R}^{m \times m}$ defines a permutation that is introduced for numerical stability [19]. In this section we first review the conventional unblocked and blocked algorithms for the LU factorization, and then describe how BDP is exploited from within them. Next, we introduce a re-organized version of the algorithm that integrates look-ahead in order to enhance performance in a nested TP+BDP execution.

All experiments in the remainder of the paper were obtained using IEEE double-precision arithmetic on an Intel Xeon E5-2603 v3 processor (6 cores at a nominal frequency 1.6 GHz). The implementations were linked with BLIS version 0.1.8 or a tailored version of this library especially developed for this work. Unless otherwise stated, BDP parallelism is extracted only from Loop 4 of the BLIS kernels.

Algorithm: $[A] := \text{LU_UNB}(A)$	Algorithm: $[A] := \text{LU_BLK}(A)$
$A \rightarrow \left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$ where A_{TL} is 0×0 while $n(A_{TL}) < n(A)$ do $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c c c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$ where α_{11} is a scalar rl1. $a_{21} := a_{21}/\alpha_{11}$ rl2. $A_{22} := A_{22} - a_{21}a_{12}^T$ $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c c c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$ endwhile	$A \rightarrow \left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$ where A_{TL} is 0×0 while $n(A_{TL}) < n(A)$ do Determine block size b $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c c c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right)$ where A_{11} is $b \times b$ RL1. $\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} := \text{LU_UNB} \left(\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} \right)$ RL2. $A_{12} := \text{TRILU}(A_{11})^{-1} A_{12}$ RL3. $A_{22} := A_{22} - A_{21} A_{12}$ $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c c c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right)$ endwhile

Figure 3: Unblocked and blocked RL algorithms for the LU factorization (left and right, respectively). In the notation, $n(\cdot)$ returns the number of columns of its argument, and $\text{TRILU}(\cdot)$ returns the strictly lower triangular part of its matrix argument, setting the diagonal entries of the result to ones.

3.1. Basic algorithms and BDP

There exist a number of algorithmic variants of the LU factorization that can accommodate partial pivoting [19]. Among these, Figure 3 (left) shows an unblocked algorithm for the so-called *right-looking* (RL) variant, expressed using the FLAME notation [23]. For simplicity, we do not include pivoting in the following description of the algorithms, though all our actual implementations, (and in particular those employed in our experimental evaluation,) integrate standard partial pivoting. The cost of computing the LU factorization of an $m \times n$ matrix, via any of the algorithms presented in this paper, is $mn^2 - n^3/3$ floating-point arithmetic operations (flops). Hereafter, we will consider square matrices of order n for which, the cost boils down to $2n^3/3$ flops. For the RL variants, the major part of these operations are concentrated in the initial iterations of the algorithm(s). For example, the first 25% iterations account for almost 58% of the flops; the first half for 87.5%; and the first 75% for more than 98%. Thus, the key to high performance mostly lies in the initial stages of the factorization.

For performance reasons, dense linear algebra libraries compute the LU factorization via a blocked algorithm that casts most computations in terms of GEMM. Figure 3 (right) presents the blocked RL algorithm. For each itera-

tion, the algorithm processes panels of b columns, where b is the algorithmic block size. The three operations in the loop body factorize the “current” panel $A_p = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix}$, via the unblocked algorithm (LU_UNB, RL1); and next update the trailing submatrix, consisting of A_{12} and A_{22} , via a triangular system solve (TRSM, RL2) and a matrix multiplication (GEMM, RL3), respectively. In practice, the block size b is chosen so that the successive invocations to the GEMM kernel deliver high FLOPS (flops per second) rates. If b is too small, the performance of GEMM will suffer, and so will that of the LU factorization. On the other hand, reducing b is appealing as this choice decreases the number of flops that are performed in terms of the panel factorization, an operation that can be expected to offer significantly lower throughput (FLOPS) than GEMM. (Concretely, provided $n \gg b$, the cost required for all panel factorizations is about $n^2b/2$ flops.) Thus, there is the tension between these two requisites.

When the target platform is a multicore processor, the conventional parallelization of the LU factorization simply relies on multi-threaded instances of TRSM and GEMM to exploit BDP only. Compared with this, the panel factorization of A_p , which lies in the critical path of the blocked RL factorization algorithm, exhibits a reduced degree of concurrency. Thus, depending on the selected block size b and certain hardware features of the target architecture (number of cores, floating-point performance, memory bandwidth, etc.), this operation may easily become a performance bottleneck; see Figure 4.

To illustrate the performance relevance of the panel factorization, Figure 5 displays a fragment of a trace corresponding to the LU factorization of a $10,000 \times 10,000$ matrix, using the blocked RL algorithm in Figure 3, with partial pivoting and “outer” block size $b = b_o = 256$. (All traces in this paper were obtained using **Extrae** version 3.3.0 [24].) The code is linked with multi-threaded versions of the BLIS kernels for GEMM and TRSM. The panel factorization (PANEL) is performed via a call to the same blocked algorithm, with “inner” block size $b = b_i = 32$, and also extracts BDP from the same two kernels. With this configuration, the panel factorization represents less than 2% of the flops performed by the algorithm. However, the trace of the first four iterations reveals that its practical cost is much higher than could be expected. (The cost of factorizing a panel relative to the cost of an iteration becomes even larger as the iteration progresses.) Here we note also the significant cost of the row permutations, which are performed via the sequential legacy code for this routine in LAPACK (LASWP). However,

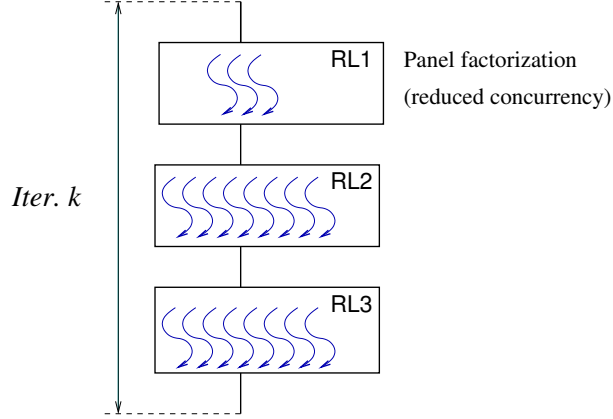


Figure 4: Exploitation of BDP in the blocked RL LU parallelization. A single thread team executes all the operations, with less active threads for RL1 due to the reduced concurrency of this kernel. In this algorithm, RL1 stands in the critical path.

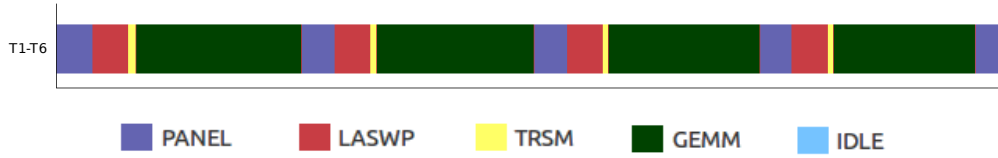


Figure 5: Execution trace of the first four iterations of the blocked RL LU factorization with partial pivoting, using 6 threads, applied to a square matrix of order 10,000, with $b_o = 256$, $b_i = 32$.

this second operation is embarrassingly parallel and its execution time can be expected to decrease linearly with the number of cores.

At this point, we note that the operations inside the loop body of the blocked algorithm in Figure 3 (right) present strict dependencies that enforce their computation in the order $RL1 \Rightarrow RL2 \Rightarrow RL3$. Therefore, there seems to be no efficient manner to formulate a TP version of the blocked algorithm in that figure.

3.2. Static look-ahead and nested TP+BDP

A strategy to tackle the hurdle represented by the panel factorization in a parallel execution consists in the introduction of look-ahead [25] into the algorithm. Concretely, during each iteration of the decomposition this technique aims to overlap the factorization of the “next” panel with the

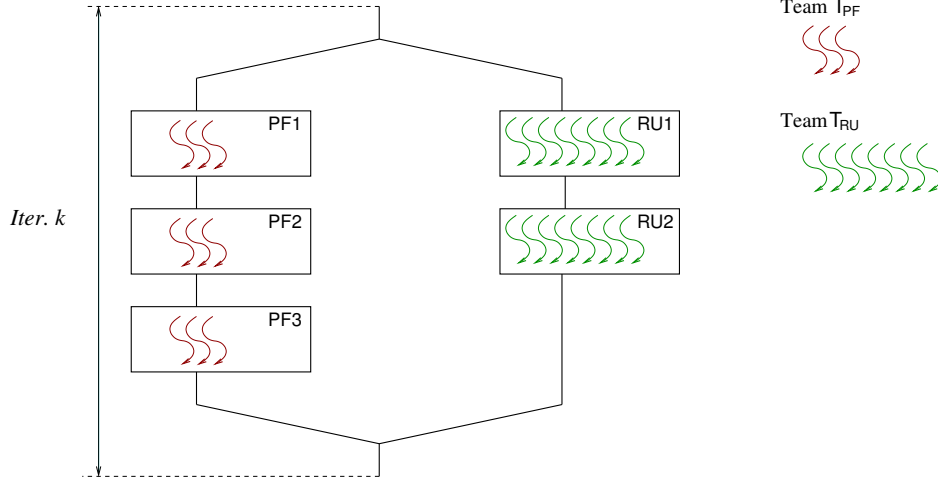


Figure 7: Exploitation of TP+BDP in the blocked RL LU parallelization with look-ahead. The execution is performed by teams T_{PF} and T_{RU} , consisting of $t_{pf} = 3$ and $t_{ru} = 8$ threads, respectively. In this algorithm, the operations on the $(k + 1)$ -th panel, including its factorization (PF3), are overlapped with the updates on the remainder of the trailing submatrix (RU1 and RU2).

of RU1 and RU2; see Figure 6. The “decoupling” of these block panels thus facilitate that, in a TP execution of an iteration of the loop body of the look-ahead version, the updates on A_{12}^P , A_{22}^P and the factorization of the latter (operations on the next panel, in T_{PF}) proceed concurrently with the updates of A_{12}^R , A_{22}^R (remainder operations, in T_{RU}), as there are no dependencies between T_{PF} and T_{RU} .

By carefully tuning the block size b and adjusting the amount of computational resources (threads) dedicated to each of the two independent tasks, T_{PF} and T_{RU} , a nested TP+BDP execution of the algorithm enhanced with this static look-ahead can partially or totally overcome the bottleneck represented by the panel factorization; see Figure 7.

Figure 8 illustrates a complete overlap of T_{RU} with T_{PF} attained by the look-ahead technique. The results in that figure correspond to a fragment of a trace obtained for the LU factorization of a $10,000 \times 10,000$ matrix, using the blocked RL algorithm in Figure 6, with partial pivoting, and outer block size $b = b_o = 256$. For this experiment, the $t = 6$ threads are partitioned into two teams: PF with $t_{pf} = 1$ thread in charge of T_{PF} , and RU with $t_{ru} = 5$ threads responsible for T_{RU} . The panel factorization (PANEL) is performed via a call to the same algorithm, with $b = b_i = 32$, and this operation proceeds

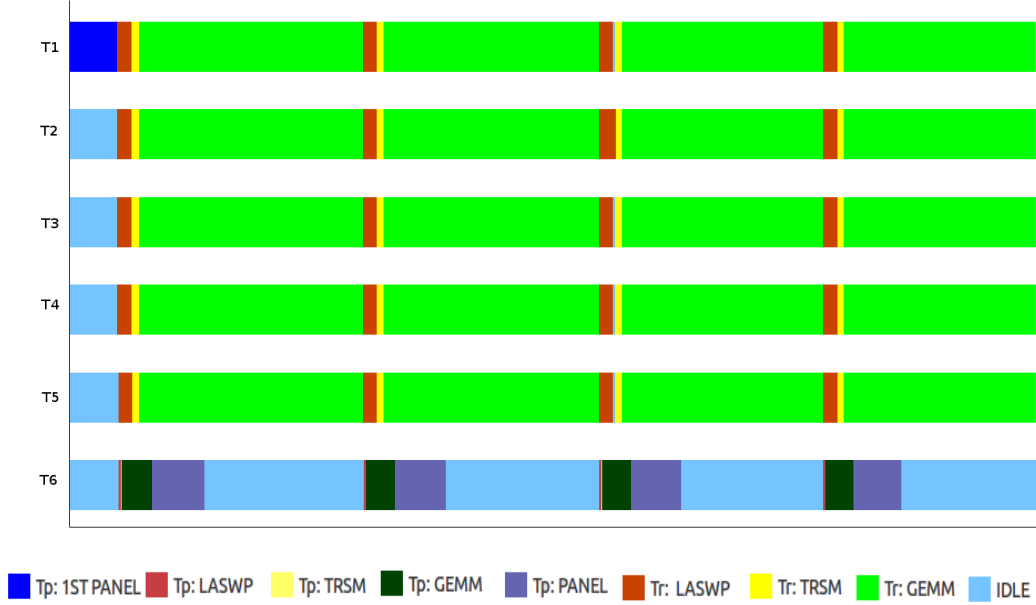


Figure 8: Execution trace of the first four iterations of the blocked RL LU factorization with partial pivoting, enhanced with look-ahead, using 6 threads, applied to a square matrix of order 10,000, with $b_o = 256$, $b_i = 32$.

sequentially (as PF consists of a single thread). The application of the row permutations is distributed between all 6 cores. As argued earlier, the net effect of the look-ahead is that the cost of the panel factorization no longer has a practical impact on the execution time of the (first four iterations of) the factorization algorithm, which is now basically determined by the cost of the remaining operations.

Given a static mapping of threads to tasks, b should balance the time spent in the two tasks as, if the operations in T_{PF} take longer than those in T_{RU} , or vice-versa, the threads in charge of the less expensive part will become idle, causing a performance degradation. This was already visible in Figure 8, which shows that, during the first four iterations, the operations in T_{PF} are considerably less expensive than the updates performed as part of the remainder T_{RU} . The complementary case, where T_{PF} requires longer than T_{RU} , is illustrated using the same configuration, for a matrix of dimension $2,000 \times 2,000$, in Figure 9. Unfortunately, as the factorization proceeds, the theoretical costs and execution times of T_{PF} and T_{RU} vary, making it difficult

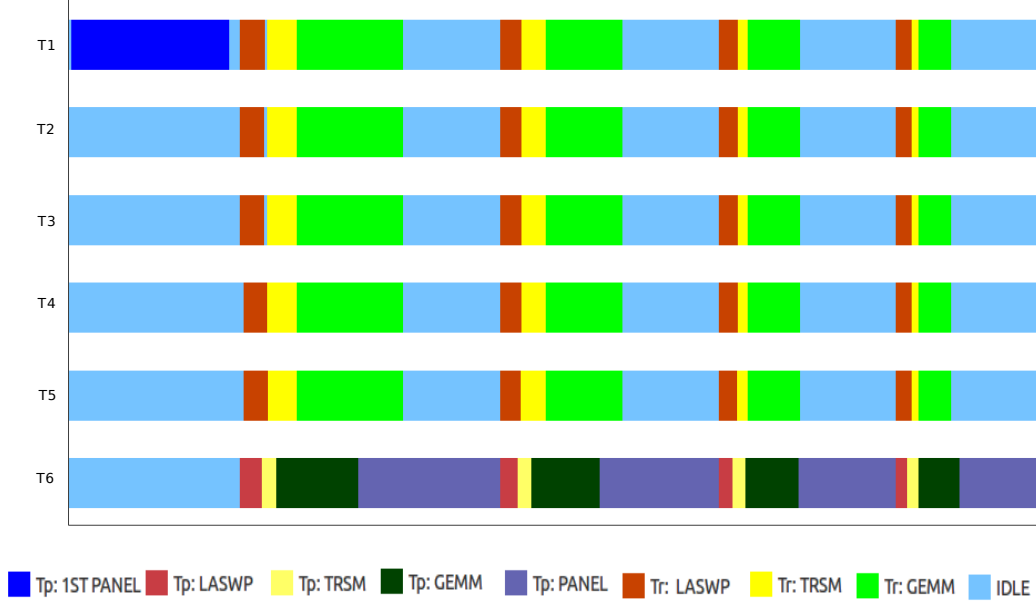


Figure 9: Execution trace of the first four iterations of the blocked RL LU factorization with partial pivoting, enhanced with look-ahead, using 6 threads, applied to a square matrix of order 2,000, with $b_o = 256$, $b_i = 32$.

to determine the optimal value of b , which will need to be adapted during the factorization process.

To close this section, note that there exist strict dependencies that serialize the operations within each task: $\text{PF1} \Rightarrow \text{PF2} \Rightarrow \text{PF3}$ and $\text{RU1} \Rightarrow \text{RU2}$. Therefore, there is no further TP in the loop-body of this re-organized version. However, the basic look-ahead mechanism of level/depth 1 described in this subsection can be refined to accommodate further levels of TP, by “advancing” to the current iteration the panel factorization of the following d iterations, in a look-ahead of level/depth d . This considerably complicates the code of the algorithm, but can be seamlessly achieved by a runtime system enhanced with priorities.

4. Advocating for Malleable Thread-Level LA Libraries

For simplicity, in the following discussion we will assume that T_{PF} and T_{RU} consist only of the panel factorization involving A_{22}^P (PF3) and the update of A_{22}^R (RU2), respectively. Furthermore, we will consider a nested TP+BDP

execution using $t = t_{\text{pf}} + t_{\text{ru}}$ threads, *initially* with a team PF of t_{pf} threads mapped to the execution of PF3 and a team RU of t_{ru} threads computing RU2.

Ideally, for the LU factorization with look-ahead, we would like to perform a *flexible sharing* of the computational resources so that, as soon as the threads in team PF complete PF3, they join team RU to help in the execution of RU2 or vice-versa. We next discuss these two cases in detail.

4.1. Worker sharing: Panel factorization less expensive than update

Our goal is to enable, at each iteration of the algorithm for the LU factorization with look-ahead, the threads in team PF that complete the panel factorization join the thread team RU working on the update. The problem is that, if the multiplication to update A_{22}^R was initiated via an invocation to a traditional GEMM, this is not possible as none of the existing high performance implementations of BLAS allow the number of threads working on a kernel that is already in execution to be modified.

4.1.1. Suboptimal solution: Static re-partitioning

A simple workaround for this problem is to split A_{22}^R into multiple column blocks, for example, $A_{22}^R \rightarrow (A_1 \mid A_2 \mid \dots \mid A_q)$, and to perform a separate call to BLAS GEMM in order to exploit BDP during the update of each block. Then, just before each invocation, the kernel’s code queries whether the execution of the panel factorization is completed and, if that is the case, executes the suboperation with the threads from both teams (or only those of RU otherwise). Unfortunately, this approach presents several drawbacks:

- Replacing a single invocation to a coarse GEMM by multiple calls to smaller GEMM may offer lower throughput because the operands passed to GEMM are smaller and/or suboptimally “shaped”. The consequence is that calling GEMM multiple times will internally incur re-packing and data movement overheads, which are more difficult to amortize because of the smaller problem dimensions.
- The burden of which loop to partition for parallelism (note that A_{22}^R could have alternatively been split by rows, or into blocks), and the granularity of this partitioning is then placed upon the programmer’s shoulders, who may lack the information that is necessary to make a good choice. For example, if the granularity is too coarse, this will have negative effect because the integration of the single thread in the

update will likely be delayed. A granularity that is too fine, on the other hand, may reduce the parallelism within the BLAS operation or result in the use of cache blocking parameters that are too small.

4.1.2. Our solution: Malleable thread-level BLAS

The alternative that we propose in this work exploits BDP inside RU2, *but allows to change the number of threads that participate in this computation even if the task is already in execution!* In other words, we view the threads as a resource pool of workers that can be shared between different tasks and reassigned to the execution of a (BLAS) kernel that is already in progress.

The key to our approach lies in the explicit exposure of the GEMM internals (and other BLAS-3 kernels) in BLIS. Concretely, assume that RU2 is computed via a single invocation to BLIS GEMM, and consider that this operation is parallelized by distributing the iteration space of Loop 4 among the threads in team PF; (see Figures 1 and 2). Then, just before Loop 4, we force the system to check if the execution of the panel factorization is completed and, based on this information, decides whether this loop is executed using either the union of the threads from both teams or only those in RU; see Figure 10.

Let us re-analyze the problems listed in Subsection 4.1.1 for the work-around that statically partitioned the update of A_{22}^R , and compare them with our solution that implicitly embeds this partitioning inside BLIS:

- The partitioning of GEMM into multiple calls to smaller matrix multiplications does not occur. Our solution performs a single call to GEMM only, so that there is no additional re-packing nor data movements. For example, in the case just discussed, B_c is already packed and re-used independently of whether t or t_{ru} threads participate in the GEMM. The buffer A_c is packed only once per iteration of Loop 3 (in parallel by both teams or only RU).
- The decision of the best partitioning/granularity is left in the hands of BLIS, which likely has more information to do a better job than the programmer.

Importantly, the partitioning happens dynamically and is transparent to the programmer.

Figure 11 validates the effect of integrating a malleable version of BLIS into the same configuration that produced the results in Figure 8. A comparison of both figures shows that, with a malleable version of BLIS, the thread

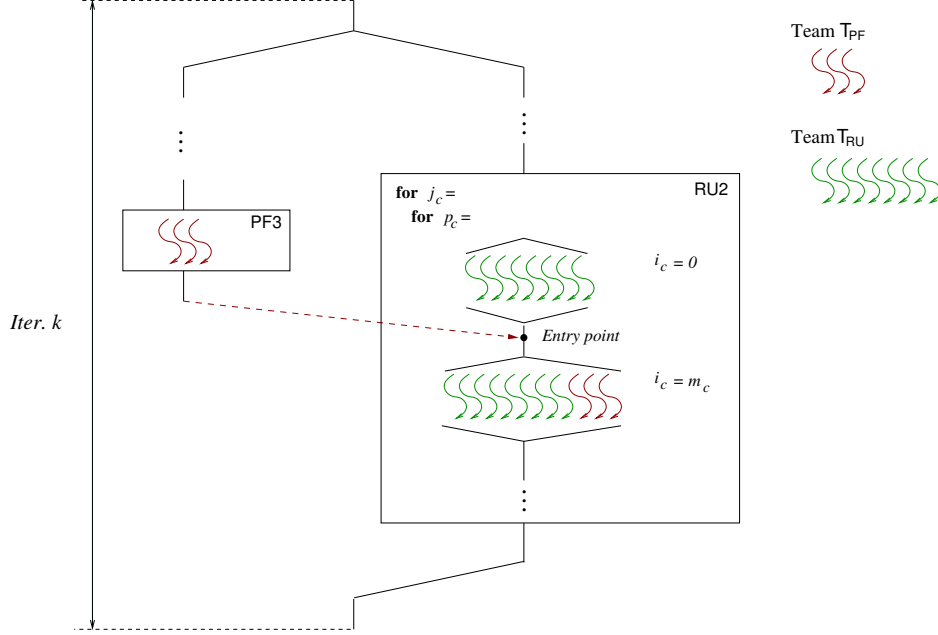


Figure 10: Exploitation of TP+BDP in the blocked RL LU parallelization with look-ahead and WS. The execution is performed by teams T_{PF} and T_{RU} , consisting of $t_{pf} = 3$ and $t_{ru} = 8$ threads, respectively. In this example, team T_{PF} completes the factorization PF3 while team T_{RU} is executing the first iteration of Loop 3 that corresponds to RU2/GEMM ($i_c = 0$). Both teams then merge and jointly continue the update of the remaining iterations of that loop ($i_c = m_c, 2m_c, \dots$). With the parallelization of GEMM Loop 4, one such “entry point” enables the merge at the beginning of each iteration of loop i_c .

executing the operations in T_{PF} , after completing this task, rapidly joins the team that computes the remainder updates, thus avoiding the idle wait.

Compared with BLIS, the same approach cannot be integrated into Goto-BLAS because the implementation of GEMM in this library only exposes the three outermost loops of Figure 1, while the remaining loops are encoded in assembly. The BLAS available as part commercial libraries is not an option either because hardware vendors offer black-box implementations which do not permit the migration of threads.

4.2. Early termination: Panel factorization more expensive than update

The analysis of this case will reveal some important insights. In order to discuss them, let us consider that, in the LU factorization with look-ahead, the panel factorization (PF3) is performed via a call to the blocked routine

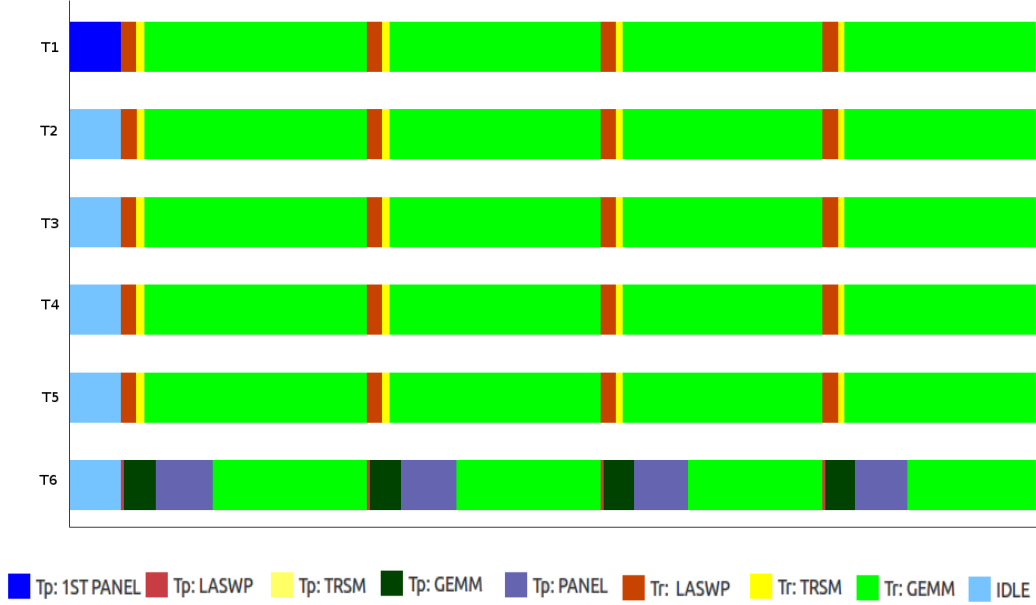


Figure 11: Execution trace of the first four iterations of the blocked RL LU factorization with partial pivoting, enhanced with look-ahead and *malleable BLIS*, using 6 threads, applied to a square matrix of order 10,000, with $b_o = 256$, $b_i = 32$.

in Figure 3 (right). We assume two blocking parameters: $b = b_o$ for the outer routine that computes the LU factorization of the complete matrix using look-ahead, and $b = b_i$ for the inner routine that factorizes each panel. (Note that, if $b_i = b_o$ or $b_i = 1$, the panel factorization is then simply done via the unblocked algorithm.) Furthermore, we will distinguish these two levels by referring to them as the *outer LU* (factorization with look-ahead) and the *inner LU* (factorization via the blocked algorithm without look-ahead). Thus, at each iteration of the outer LU, a panel of b_o columns is factorized via a call to LU_BLK (inner LU), and this second decomposition proceeds to factorize the panel using a blocked algorithm with block size b_i ; see Figure 12.

From Figure 3 (right), the loop body for the inner LU consists of a call to the unblocked version of the algorithm (RL1), followed by the invocations to TRSM and GEMM that update A_{12} and A_{22} , respectively (RL2 and RL3). Now, let us assume that the update RU2 by the thread team RU is completed while the threads of team PF are in the middle of the computations corresponding to an iteration of the loop body of the inner LU. Then, provided the versions

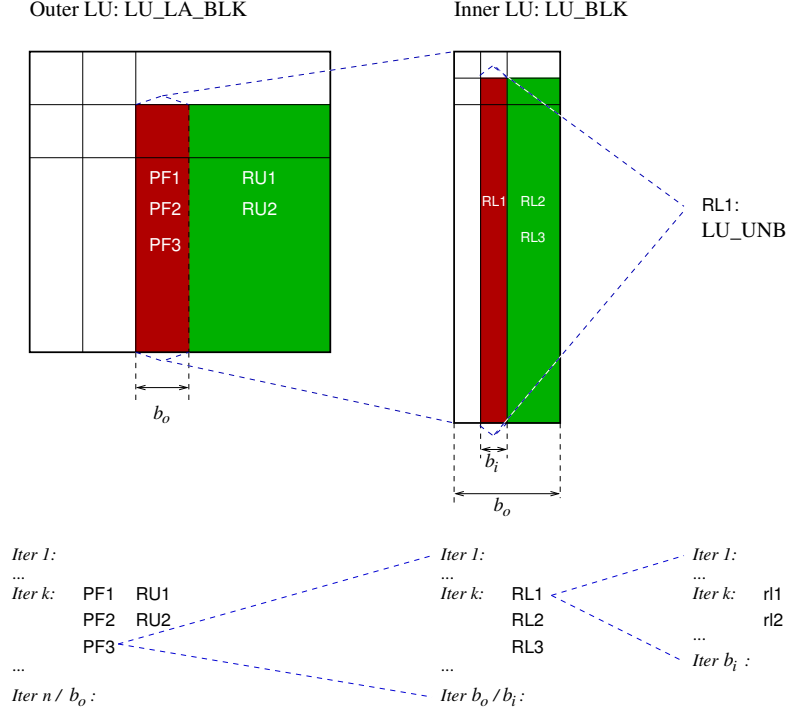


Figure 12: Outer vs inner LU and use of algorithmic block sizes.

of the BDP versions TRSM and GEMM kernels that are invoked from the inner LU are malleable (see subsection 4.1.2), inside them the system will perform the actions that are necessary to integrate the thread team RU, which is now idle, into the corresponding (and subsequent) computation(s). Unfortunately, the updates in the loop body of the inner LU involve small-grained computations (A_{12} and A_{22} have at most $b_o - b_i$ columns, decreasing by b_i columns at each iteration), and little parallel performance can be expected from it especially because of partial pivoting.

In order to deal with this scenario, a different option is to force the inner LU to stop at the end of the current iteration, to then rapidly proceed to the next iteration of the outer LU. We refer to this strategy as the *early termination* (ET). In order to do this though, the transformations computed to factorize the current inner-panel must be propagated first to the remaining columns outside this panel, introducing a certain delay in this version of the ET strategy.

A third possibility is to rely on a left-looking (LL) version of the LU

factorization for the inner LU, as discussed next. The blocked LL algorithm for the LU factorization differs from the blocked RL variant (see the algorithm in the right-hand side of Figure 3) in the operations performed inside the loop-body, which are replaced by

$$\begin{aligned}
\text{LL1.} \quad A_{01} &:= \text{TRILU}(A_{00})^{-1} A_{01} \\
\text{LL2.} \quad \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} &:= \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} - \begin{bmatrix} A_{10} \\ A_{20} \end{bmatrix} A_{01} \\
\text{LL3.} \quad \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} &:= \text{LU_UNB} \left(\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} \right)
\end{aligned}$$

Thus, at the end of a certain iteration, this variant has only updated the current column of the inner-panel and those to its left. In other words, no transformations are propagated beyond that point (i.e., to the right of the current column/inner-panel), and ET can be implemented in a straightforward manner, with no delay compared with an inner LU factorization via the RL variant.

A definitive advantage of the LL variant compared with its RL counterpart is that the former implements a lazy algorithm, which delays the operations towards the end of the panel factorization, while the second corresponds to an eager algorithm that advances as much computations as possible to the initial iterations. Therefore, in case the panel factorization has to be stopped early, it is more likely that the LL variant has progressed in the factorization further.³ The appealing consequence is that this enables the use of larger block sizes for the following updates in the LL variant.

From an implementation point of view, the synchronization between the two teams of threads is easy to handle. For example, at the beginning of each iteration of the outer LU, a boolean flag is set to indicate that the remainder update is incomplete. The thread team RU then changes this value as soon as this task is complete. In the mean time, the flag is queried by the thread team PF, at the end of every iteration of the inner LU, aborting its execution when a change is detected. With this operation mode, there is no need to protect the flag from race conditions. This solution also provides an adaptive

³Consider the factorization of an $m \times n$ matrix that is stopped at iteration $k < n$. The LL algorithm will have performed $m^2k - m^3/3$ flops at that point while, for the RL algorithm, the flop count raises to that of the LL algorithm plus $2(n-k)(mk - k^2/2)$.

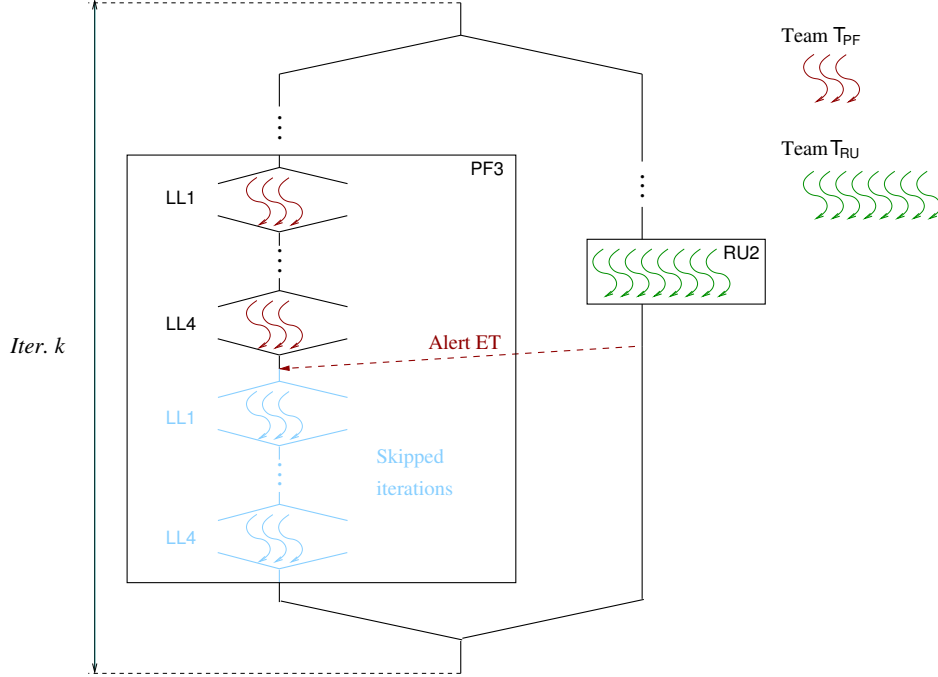


Figure 13: Exploitation of TP+BDP in the blocked RL LU parallelization with look-ahead and ET. The execution is performed by teams T_{PF} and T_{RU} , consisting of $t_{pf} = 3$ and $t_{ru} = 8$ threads, respectively. In this example, team T_{RU} completes the update RU2 while team T_{PF} is executing an iteration of the panel factorization PF3. T_{RU} then notifies of this event to T_{PF} , which then skips the remaining iterations of the loop that processes the panel.

(automatic) configuration of the block size as, if chosen too large, it will be adjusted for the current (and, possibly, subsequent) iterations by the early termination of the inner LU. The process is illustrated in Figure 13.

4.3. Relation to adaptive look-ahead via a runtime

Compared with our approach, which only applies look-ahead at one level, a TP execution that relies on a run-time for adaptive-depth look-ahead exposes a higher degree of parallelism from “future iterations”, which can amortize the cost of the panel factorization over a considerably larger number of flops. This can be beneficial for architectures with a large number of cores, but can be partially compensated by increasing the number of threads dedicated to the panel factorization, combined with a careful fine-grain exploitation of the concurrency [26], in our approach. On the other hand,

adaptive-depth look-ahead via a runtime suffers from re-packing and data movement overheads due to multiple calls to GEMM. Moreover, it couples the algorithmic block size that fixes the granularity of the tasks to that of the suboperands in GEMM. Finally, runtime-based solutions rarely exploit nested TP+BDP parallelism and, even if they do so, taking advantage of a malleable thread-level BLAS from within them may be difficult.

5. Experimental Evaluation

In this section we analyze in detail the performance behavior of several multi-threaded implementations of the algorithms for the LU factorization:

- LU: Blocked RL (Figure 3). This code only exploits BDP, via calls to the (non-malleable) multi-threaded BLIS (version 0.1.8).
- Variants enhanced with look-ahead (Figure 6). The following three implementations take advantage of nested TP+BDP, with 1 thread dedicated to the operations on the panel (team T_{PF}) and $t - 1$ to the remainder updates (team T_{RU}).
 - LU_LA (subsection 3.2): Blocked RL with look-ahead.
 - LU_MB (subsection 4.1.2): Blocked RL with look-ahead and malleable BLIS.
 - LU_ET (subsection 4.2): Blocked RL with look-ahead, malleable BLIS, and early termination of the panel factorization.
- LU_OS: Blocked RL with adaptive look-ahead extracted via the OmpSs runtime (version 16.06). LU_OS decomposes the factorization into a large collection of tasks connected via data dependencies, and then exploits TP only, via calls to a sequential instance of BLIS. In more detail, the OmpSs parallel version divides the matrix into a collection of panels of fixed width b_o . All operations performed during an iteration of the algorithm on the same panel (row permutation, triangular system solve, matrix multiplication and, possibly, panel factorization) are then part of the same task. This implementation includes priorities to advance the schedule of tasks involving panel factorizations.

All codes include standard partial pivoting and compute the same factorization. Also, all solutions perform the panel factorization via the blocked

RL algorithm, except for LU_ET and LU_OS, which employ the blocked LL variant. The performance differences between the LL and RL variants, when applied solely to the panel factorization, were small. Nonetheless, for LU_ET, employing the LL variant improves the ET mechanism and unleashes a faster execution of the global factorization. For LU_OS we integrated the LL variant as well to favor a fair comparison between this implementation and our LU_ET. The block size is fixed to b_o during the complete iteration in all cases, except for LU_ET which initially employs b_o , but then adjusts this value during the factorization as part of the ET mechanism.

In the experiments, we considered the factorization of square matrices, with random entries uniformly distributed in $(0, 1)$, and dimension $n=500$ to 12,000 in steps of 500. The block size for the outer LU was tested for values $b_o=32$ to 512 in steps of 32. The block size for the inner LU was evaluated for $b_i=16$ and 32. We employed one thread per core (i.e., $t = 6$) in all executions.

5.1. Optimal block size

The performance of the blocked LU algorithms is strongly influenced by the outer block size b_o . As discussed in subsection 3.1, this parameter should balance two criteria:

- Deliver high performance for the GEMM kernel. Concretely, in the algorithms in Figures 3 and 6, a value of b_o that is too small turns A_{21} and A_{12}/A_{12}^R into narrow column and row panels respectively, transforming the matrix multiplication involving these blocks (RL3 in Figure 3 and RU2 in Figure 6) into a memory-bound kernel that will generally deliver low performance. In the following, we will refer to a GEMM (1) with dimensions $m \approx n \gg k$, as a panel-panel multiplication (GEPP) [8]. Note that, for the GEPP arising in the LU factorizations, $k = b_o$.
- Reduce the amount of operations performed in the panel factorization (about $n^2 b_o / 2$ flops, provided $n \gg b_o$), in order to avoid the negative impact of this intrinsically sequential stage.

Figure 14 sheds further light on the roles played by these two factors. The plot in the left-hand side reports the performance of GEPP, in terms of GFLOPS (billions of FLOPS), showing that the implementation of this kernel

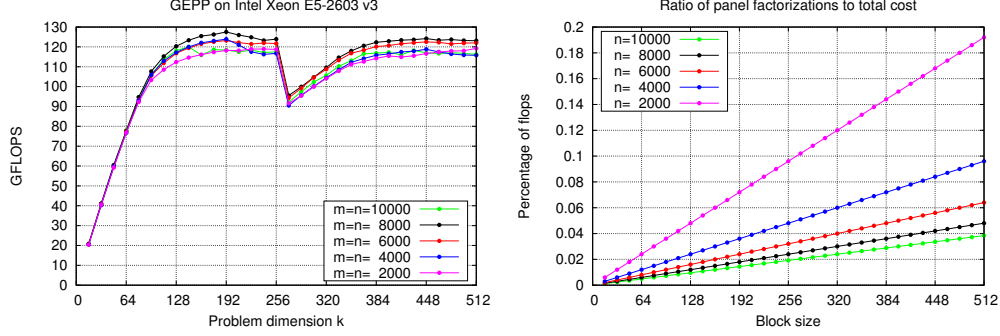


Figure 14: GFLOPS attained with GEPP (left) and ratio of flops performed in the panel factorizations normalized to the total cost (right).

in BLIS achieves an asymptotic performance peak for $k(=b_o)$ around 144.⁴ The right-hand side plot reports the ratio of flops performed in the panel factorizations with respect to those of the LU factorization.

The combined effect of these criteria seems to point in the direction of choosing the smallest b_o that attains the asymptotic GFLOPS rate for GEPP. However, Figure 15 illustrates the experimental optimal block size b_o for the distinct LU factorization algorithms, exposing that this is not the case. We next discuss the behavior for LU, LU_LA and LU_MB, which show different trends. (LU_ET and LU_OS are analyzed latter.) In particular, LU benefits from the use of larger values of b_o than the other two codes for all problem dimensions. The reason is that a large block size operates on wide panels, which turns their factorization into a BLAS-3 operation with a mild degree of parallelism, and reduces the impact of this computation on the critical path of the factorization. LU_LA exhibits a similar behavior for large problems, but favors smaller block sizes for small to moderate problems. The reason is that, for LU_LA, it is important to balance the panel factorization (T_{PF}) and remainder update (T_{RU}) so that their execution approximately requires the same time.

Compared with the previous two implementations, LU_MB promotes the use of small block sizes, up to $b_o=192$, for the largest problems. (Interestingly, this corresponds to the optimal value of k for GEPP.) One reason

⁴The performance drop observed for k slightly above 256 is due to the optimal value of k_c being equal to that number in this architecture.

for this behavior is that, when the malleable version of BLIS is integrated into LU_MB, the practical costs of the two branches/tasks do not need to be balanced. Let us elaborate this case further, by considering the effect of reducing the block size, for example, from b_o to $b'_o = b_o/2$. For simplicity, in the following discussion we will use approximations for the block dimensions and their costs; furthermore, we will assume that $n \gg b_o$. The first and most straight-forward consequence of halving the block size is that the number of iterations is doubled. Inside each iteration with the original block size b_o , the loop body invokes, among others kernels, a GEPP of dimensions $m \times (m - b_o) \times b_o$ (with m the number of rows in the trailing submatrix A_{22}^R), for a cost of $2m^2b_o$ flops; in parallel, the factorization involves a panel of dimension $m \times b_o$, for a cost of $mb_o^2 - b_o^3/3 \approx mb_o^2$ flops. When the block size is halved to b'_o , the same work is basically computed in two consecutive iterations. However, this reduces the amount of flops performed in terms of panel factorizations to about $2m(b'_o)^2 = mb_o^2/2$ while it has a minor impact on the number of flops that are cast as GEPP (two of these products, at a cost of $2m^2b'_o = 2m^2b_o/2$ flops each). The conclusion is that, by reducing the block size, we decrease the time that the single thread spends in the panel factorization T_{PF} , favoring its rapid merge with the thread team that performs the remainder update T_{RU} . Thus, in case the execution time of the LU is dominated by T_{RU} , adding one more thread to perform this task (in this scenario, in the critical path) as soon as possible will reduce the global execution time of the algorithm.

5.2. Performance comparison of the variants with static look-ahead

The previous analysis on the effect of the block size exposes that choosing the optimal block size is a difficult task. Either we need a model that can accurately predict the performance of each building block appearing in the LU factorization, or we perform an extensive experimental analysis to select the best value. The problem is even more complex if we consider that, in practice, an optimal selection would have to vary the block size as the factorization progresses. Concretely, for the factorization of a square matrix of order n via a blocked algorithm, the problem is decomposed into multiple subproblems that involve the factorization of matrices of orders $n - b_o$, $n - 2 \cdot b_o$, $n - 3 \cdot b_o$, etc. From Figure 15, it is clear that the optimal value of b_o will be different for several of these subproblems. In the end, the value that we show in Figure 15 for each problem has to be considered as a compromise

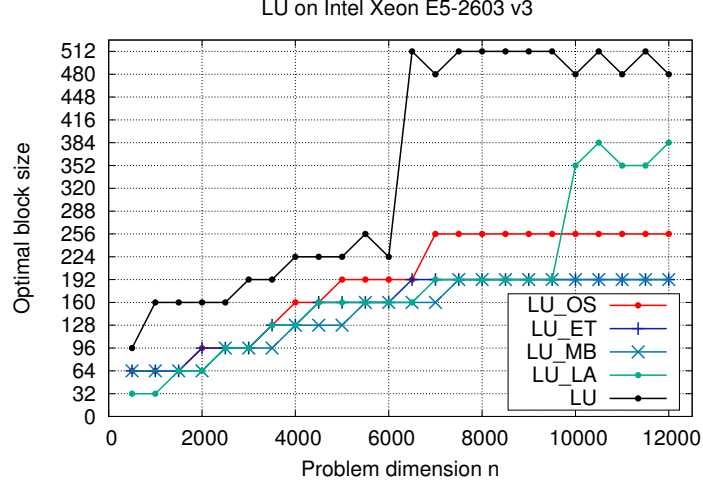


Figure 15: Optimal block size of the blocked RL algorithms for the LU factorization.

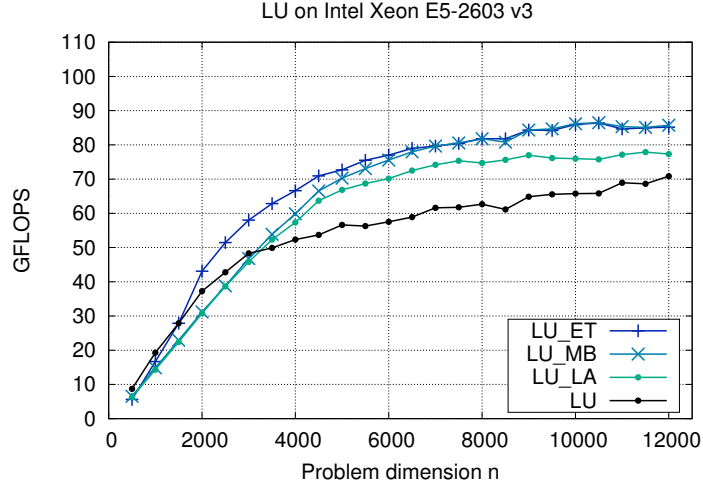


Figure 16: Performance comparison of the blocked RL algorithms for the LU factorization (except LU_OS) with a fixed block size $b_o = 256$.

that attains fair performance for a wide range of the subproblems appearing in that case.

Figure 16 reports the GFLOPS rates attained by the distinct implementations to compute the plain LU factorization and the variants equipped with static look-ahead (i.e., all except LU_OS), using $b_o = 256$ as a compromise

value for all of them. Although this value is optimal for only a few cases, the purpose of this experiment is to show the improvements attained by gradually introducing the techniques enhancing look-ahead. The figure reveals some relevant trends:

- Except for the smallest problems, integrating the look-ahead techniques clearly improves the performance of the plain LU factorization implemented in LU.
- The version with malleable BLAS (LU_MB) improves the performance of the basic version of look-ahead (LU_LA) for the larger problems. This is a consequence of the cost of the panel factorization relative to that of the global factorization. Concretely, for fixed b_o , as the problem size grows, the global flop-cost varies cubically in n , as $2n^3/3$, while the flop-cost of the panel factorizations grows quadratically, with $n^2b_o/2$. Thus, we can expect that, for large n , the remainder update T_{RU} becomes more expensive than the panel factorization T_{PF} . This represents the actual scenario that was targeted by the variant with malleable BLIS.
- The version that combines the malleable BLAS with ET (LU_ET) delivers the same performance of LU_MB for large problems, but outperforms all other variants with static look-ahead for the smaller problems. Again, this could be expected by considering the relative cost of the panel factorization for small n .

5.3. Performance comparison with OmpSs

We conclude the experimental analysis by providing a comparison of the best variant with static look-ahead, LU_ET, with the implementation that extracts parallelism via the OmpSs runtime, LU_OS. In this last experiment we depart from the previous case, performing an extensive evaluation in order to report the performance for the optimal block size for each problem dimension and algorithm. (See Figure 15 for the actual optimal values employed in the experiment.) For LU_OS, we select a value for b_o that is then fixed for the complete factorization. As this variant overlaps the execution of tasks from different iterations in time, it is difficult to vary the block size as the factorization progresses. For LU_ET, the selected value of b_o only applies to the first factorization. After that, the ET mechanism automatically adjusts this value during the iteration.

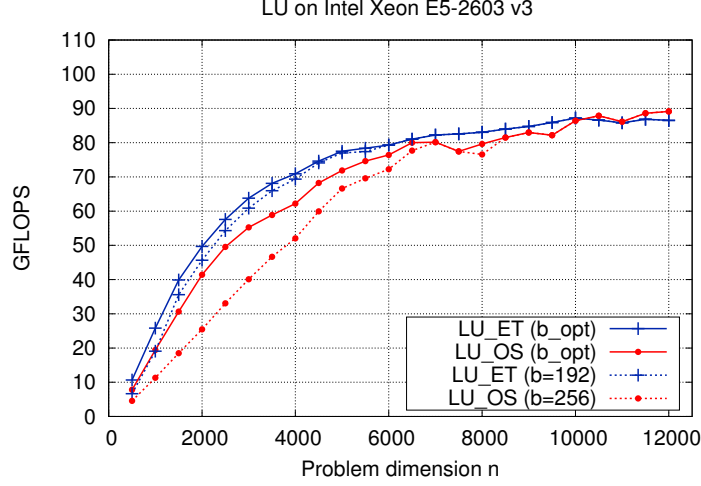


Figure 17: Performance comparison between the OmpSs implementation and the blocked RL algorithm for the LU factorization with look-ahead, malleable BLIS and ET. Two configurations are chosen for each algorithm: optimal block size for each problem size; and fixed block sizes $b_o = 192$ for LU_ET and $b_o = 256$ for LU_OS.

Figure 17 shows the results for this comparison in the lines labelled as “(b_opt)”. LU_ET is very competitive, clearly outperforming the runtime-based solution for most problems and offering competitive performance for the largest three.

Manually tuning the block size to each problem dimension is in general impractical. For this reason, the figure also shows the performance curves when the block size is fixed to $b_o = 192$ for LU_ET and $b_o = 256$ for LU_OS. These values were selected because they offered high performance for a wide range of problem dimensions (especially, the largest ones; see Figure 15). Interestingly, the performance lines corresponding to this configuration, labelled with “(b=192)”/“(b=256)”, show that choosing a suboptimal value for b_o has a minor impact on the performance of our solution LU_ET, because the ET mechanism adjust this value on-the-fly (for the smaller problem sizes). Compared with this, the negative effect of a suboptimal selection on LU_OS is clearly more visible.

A comparison with other parallel versions of the LU factorization with partial pivoting is possible, but we do not expect the results change the message of our paper. In particular, Intel MKL includes a highly-tuned routine for this factorization that relies in their own implementation of the BLAS

and some type of look-ahead. Therefore, whether the advantages of one implementation over the other come simply from the use of a different version of the BLAS, or from the positive effects of our WS and ET mechanism, will be really difficult to infer. The PLASMA library [17] also provides a routine for the LU factorization with partial pivoting supported by a runtime that implements dynamic look-ahead. The techniques integrated in PLASMA’s routine are not different from those in the OmpSs implementation evaluated in our paper. Therefore, when linked with BLIS, we do not expect a different behaviour between PLASMA’s routine and LU_OS.

6. Concluding Remarks and Future Work

We have introduced WS and ET as two novel techniques to avoid workload imbalance during the execution of matrix factorizations, enhanced with look-ahead, for the solution of linear systems. The WS mechanism especially benefits from the adoption of a malleable thread-level instance of BLIS, which allows the thread team in charge of the panel factorization, upon completion of this task, to be reallocated to the execution of the trailing update. The ET mechanism tackles the opposite situation, with a panel factorization that is costlier than the trailing update. In such scenario, the team that performed the update communicates to the second team that it should terminate the panel factorization, advancing the factorization process into the next iteration.

Our results on an Intel Xeon E5-2603 v3 show the performance benefits of our version enhanced with malleable BLIS and ET compared with a plain LU factorization as well as a version with look-ahead. The experiments also report competitive performance compared with an LU factorization that is parallelized by means of a sophisticated runtime, such as OmpSs, that introduces look-ahead of dynamic (variable) depth. Compared with the OmpSs solution, our approach offers higher performance for most problem dimensions, seamlessly tunes the algorithmic block size, and features a considerably smaller memory footprint as it does not require a sophisticated runtime support.

To conclude, our paper does not intend to propose an alternative to runtime-based solutions. Instead, the message implicitly carried in our experiments aims to emphasize the benefits of malleable thread-level libraries, which we expect to be crucial in order to exploit the massive thread parallelism of future architectures. This work opens a plethora of interesting

questions for future research. In particular, how to generalize the ideas to a multi-task scenario, what kind of interfaces may ease thread-level malleability, and what kind of support is necessary in the runtime for this purpose.

Acknowledgements

We thank the other members of the FLAME team for their support. This research was partially sponsored by projects TIN2014-53495-R and TIN2015-65316-P of the Spanish *Ministerio de Economía y Competitividad*, the H2020 EU FETHPC Project 671602 “INTERTWinE”, by project 2014-SGR-1051 from the Generalitat de Catalunya, and NSF grant ACI-1550493.

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation (NSF).

References

- [1] C. L. Lawson, R. J. Hanson, D. R. Kincaid, F. T. Krogh, Basic linear algebra subprograms for Fortran usage, ACM Trans. Math. Soft. 5 (3) (1979) 308–323.
- [2] J. J. Dongarra, J. Du Croz, S. Hammarling, R. J. Hanson, An extended set of FORTRAN basic linear algebra subprograms, ACM Trans. Math. Softw. 14 (1) (1988) 1–17.
- [3] J. J. Dongarra, J. Du Croz, S. Hammarling, I. Duff, A set of level 3 basic linear algebra subprograms, ACM Trans. Math. Softw. 16 (1) (1990) 1–17.
- [4] Intel, Math Kernel Library, <https://software.intel.com/en-us/intel-mkl> (2015).
- [5] AMD, AMD Core Math Library, <http://developer.amd.com/tools-and-sdks/cpu-development/amd-core-math-library-acml/> (2015).
- [6] IBM, Engineering and Scientific Subroutine Library, <http://www-03.ibm.com/systems/power/software/essl/> (2015).
- [7] NVIDIA, cublas, <https://developer.nvidia.com/cublas> (2016).

- [8] K. Goto, R. A. v. d. Geijn, Anatomy of high-performance matrix multiplication, *ACM Trans. Math. Softw.* 34 (3) (2008) 12:1–12:25.
doi:10.1145/1356052.1356053.
URL <http://doi.acm.org/10.1145/1356052.1356053>
- [9] K. Goto, R. van de Geijn, High performance implementation of the level-3 BLAS, *ACM Transactions on Mathematical Software* 35 (1) (2008) 4:1–4:14.
URL <http://doi.acm.org/10.1145/1377603.1377607>
- [10] <http://www.openblas.net> (2015).
- [11] R. C. Whaley, J. J. Dongarra, Automatically tuned linear algebra software, in: *Proceedings of SC'98*, 1998.
- [12] F. G. Van Zee, R. A. van de Geijn, BLIS: A framework for rapidly instantiating BLAS functionality, *ACM Trans. Math. Softw.* 41 (3) (2015) 14:1–14:33.
- [13] E. Anderson, Z. Bai, L. S. Blackford, J. Demmesl, J. J. Dongarra, J. D. Croz, S. Hammarling, A. Greenbaum, A. McKenney, D. C. Sorensen, *LAPACK Users' guide*, 3rd Edition, SIAM, 1999.
- [14] F. G. V. Zee, *libflame: The Complete Reference*, www.lulu.com, 2009.
- [15] OmpSs project home page, <http://pm.bsc.es/ompss>.
- [16] StarPU project, <http://runtime.bordeaux.inria.fr/StarPU/>.
- [17] PLASMA project home page, <http://icl.cs.utk.edu/plasma>.
- [18] FLAME project home page, <http://www.cs.utexas.edu/users/flame/>.
- [19] G. H. Golub, C. F. V. Loan, *Matrix Computations*, 3rd Edition, The Johns Hopkins University Press, Baltimore, 1996.
- [20] F. G. V. Zee, T. M. Smith, B. Marker, T. M. Low, R. A. V. D. Geijn, F. D. Igual, M. Smelyanskiy, X. Zhang, M. Kistler, V. Austel, J. A. Gunnels, L. Killough, The BLIS framework: Experiments in portability, *ACM Trans. Math. Softw.* 42 (2) (2016) 12:1–12:19.
doi:10.1145/2755561.
URL <http://doi.acm.org/10.1145/2755561>

- [21] T. M. Smith, R. van de Geijn, M. Smelyanskiy, J. R. Hammond, F. G. Van Zee, Anatomy of high-performance many-threaded matrix multiplication, in: Proc. IEEE 28th Int. Parallel and Distributed Processing Symp., IPDPS'14, 2014, pp. 1049–1059.
- [22] S. Catalán, F. D. Igual, R. Mayo, R. Rodríguez-Sánchez, E. S. Quintana-Ortí, Architecture-aware configuration and scheduling of matrix multiplication on asymmetric multicore processors, *Cluster Computing* 19 (3) (2016) 1037–1051. doi:10.1007/s10586-016-0611-8.
URL <http://dx.doi.org/10.1007/s10586-016-0611-8>
- [23] J. A. Gunnels, F. G. Gustavson, G. M. Henry, R. A. van de Geijn, FLAME: Formal linear algebra methods environment, *ACM Trans. Math. Softw.* 27 (4) (2001) 422–455.
- [24] B. S. Center, Extrae user guide manual for version 2.5.1 (2016).
- [25] P. Strazdins, A comparison of lookahead and algorithmic blocking techniques for parallel matrix factorization, Tech. Rep. TR-CS-98-07, Department of Computer Science, The Australian National University, Canberra 0200 ACT, Australia (1998).
- [26] A. M. Castaldo, R. C. Whaley, S. Samuel, Scaling LAPACK panel operations using parallel cache assignment, *ACM Trans. Math. Soft.* 39 (4) (2013) 22:1–22:30.